

WHAT IS CLAIMED IS:

1. A method, comprising:

5 a load balancer receiving a request;

the load balancer selecting a node to handle the request from among a plurality of nodes associated with the load balancer and not known by the load balancer to be inactive;

10

the load balancer determining if the selected node is able to service the request;

if the selected node is determined to be unable to service the request, the load balancer selecting another node to handle the request from among the plurality of nodes associated with the load balancer and not known by the load balancer to be inactive.

15

2. The method as recited in claim 1, wherein the load balancer is one load balancer among a plurality of load balancers in a load balancer hierarchy.

20

3. The method as recited in claim 2, wherein the plurality of nodes associated with the load balancer are load balancers in a lower-level of the load balancer hierarchy.

4. The method as recited in claim 2, wherein the load balancer is associated with a higher-level load balancer in the load balancer hierarchy, and wherein said receiving a request comprises receiving the request from the higher-level load balancer.

25

5. The method as recited in claim 4, further comprising, if the selected node is determined to be unable to service the request and if no other nodes from among the plurality of nodes associated with the load balancer are not known by the load balancer to

30

be inactive, the load balancer sending a message to the higher-level load balancer to disable the load balancer from receiving further requests.

5 6. The method as recited in claim 5, further comprising, upon receiving said message, the higher-level load balancer marking the load balancer as inactive.

10 7. The method as recited in claim 5, further comprising, upon receiving said message, the higher-level load balancer re-load-balancing requests pending with the load balancer among other load balancers associated with the higher-level load balancer.

 8. The method as recited in claim 1, wherein said determining if the selected node is able to service the request comprises the load balancer actively probing the plurality of nodes associated with the load balancer.

15 9. The method as recited in claim 8, further comprising the load balancer periodically performing said actively probing.

20 10. The method as recited in claim 8, further comprising, if one of the plurality of nodes associated with the load balancer does not respond to said active probing within a timeout period, the load balancer marking that node as inactive.

25 11. The method as recited in claim 10, wherein the load balancer marking that node as inactive comprises re-load-balancing requests pending with that node among the plurality of nodes associated with the load balancer and not known by the load balancer to be inactive.

 12. The method as recited in claim 10, wherein the load balancer marking that node as inactive comprises, if no other nodes from among the plurality of nodes associated with the load balancer are not known by the load balancer to be inactive, the

load balancer sending a message to the higher-level load balancer to disable the load balancer from receiving further requests.

13. The method as recited in claim 1, further comprising:

5

the load balancer sending the request to the selected node;

wherein said determining if the selected node is able to service the request comprises the load balancer determining if the selected node fails to respond to the request within a timeout period.

10

14. The method as recited in claim 13, further comprising, if the selected node fails to respond to the request within the timeout period, the load balancer marking the selected node as inactive.

15

15. The method as recited in claim 14, wherein the load balancer marking the selected node as inactive comprises, if no other nodes from among the plurality of nodes associated with the load balancer are not known by the load balancer to be inactive, the load balancer sending a message to the higher-level load balancer to disable the load balancer from receiving further requests.

20

16. The method as recited in claim 14, wherein the load balancer marking the selected node as inactive comprises re-load-balancing requests pending with the selected node among the plurality of nodes associated with the load balancer and not known by the load balancer to be inactive.

25

17. The method as recited in claim 1, further comprising:

after said selecting the node, the load balancer sending a dummy request to the selected node;

30

wherein said determining if the selected node is able to service the request comprises the load balancer determining if the selected node fails to respond to the dummy request within a timeout period.

5

18. The method as recited in claim 17, further comprising if the selected node fails to respond to the dummy request within the timeout period, the load balancer marking the selected node as inactive.

10

19. The method as recited in claim 18, wherein the load balancer marking the selected node as inactive comprises, if no other nodes from among the plurality of nodes associated with the load balancer are not known by the load balancer to be inactive, the load balancer sending a message to the higher-level load balancer to disable the load balancer from receiving further requests.

15

20. The method as recited in claim 18, wherein the load balancer marking the selected node as inactive comprises re-load-balancing requests pending with the selected node among the plurality of nodes associated with the load balancer and not known by the load balancer to be inactive.

20

21. The method as recited in claim 17, further comprising, if the selected node responds to the dummy request within the timeout period, the load balancer sending the request to the selected node.

25

22. The method as recited in claim 21, wherein said determining if the selected node is able to service the request further comprises the load balancer determining if the selected node fails to respond to the request within a timeout period.

23. The method as recited in claim 1, wherein said determining if the selected node is able to service the request comprises the load balancer receiving a message from the selected node indicating that the selected node is disabled.

5 24. The method as recited in claim 23, further comprising, upon receiving said message, the load balancer marking the selected node as inactive.

25. The method as recited in claim 24, further comprising, upon receiving said message, the load balancer re-load-balancing requests pending with the selected node
10 among the plurality of nodes associated with the load balancer and not known by the load balancer to be inactive.

26. A system, comprising:

15 a plurality of nodes;

a load balancer associated with said plurality of nodes, wherein the load balancer is configured to:

20 receive a request;

select a node to handle the request from the plurality of nodes, wherein the plurality of nodes are not known by the load balancer to be inactive;

25 determine if the selected node is able to service the request;

select another node to handle the request from among the plurality of nodes not known by the load balancer to be inactive if the selected
30 node is determined to be unable to service the request.

27. The system of claim 26 further comprising a load balancer hierarchy, wherein the load balancer is one load balancer among a plurality of load balancers in the load balancer hierarchy.

5

28. The system of claim 27, wherein the plurality of nodes are load balancers in a lower-level of the load balancer hierarchy.

29. The system of claim 27, wherein the load balancer is associated with a
10 higher-level load balancer in the load balancer hierarchy, and wherein the load balancer is configured to receive the request from the higher-level load balancer.

30. The system of claim 29 wherein the load balancer is further configured to
send a message to the higher-level load balancer to disable the load balancer from
15 receiving further requests if the selected node is determined to be unable to service the request and if no other nodes from among the plurality of nodes associated with the load balancer are not known by the load balancer to be inactive.

31. The system of claim 30 wherein the higher-level load balancer is
20 configured to mark the load balancer as inactive upon receiving said message.

32. The system of claim 30 wherein the higher-level load balancer is
configured to re-load-balance requests pending with the load balancer among other load
balancers associated with the higher-level load balancer upon receiving said message.

25

33. The system of claim 26, wherein to determine if the selected node is able
to service the request, the load balancer is configured to actively probe the plurality of
nodes associated with the load balancer.

34. The system of claim 33, wherein the load balancer is configured to periodically actively probe the plurality of nodes associated with the load balancer.

35. The system of claim 33 wherein the load balancer is configured to mark
5 one of the plurality of nodes associated with the load balancer as inactive if that node does not respond to the active probe within a timeout period.

36. The system of claim 35, wherein the load balancer is configured to re-load-balance requests pending with the inactive node among the plurality of nodes
10 associated with the load balancer and not known by the load balancer to be inactive.

37. The system of claim 35, wherein, if no other nodes from among the plurality of nodes associated with the load balancer are not known by the load balancer to be inactive, the load balancer is configured to send a message to the higher-level load
15 balancer to disable the load balancer from receiving further requests.

38. The system of claim 26 wherein the load balancer is further configured to send the request to the selected node; and to determine if the selected node is able to service the request, the load balancer is configured to determine if the selected node fails
20 to respond to the request within a timeout period.

39. The system of claim 38 wherein the load balancer is configured to mark the selected node as inactive if the selected node fails to respond to the request within the timeout period.

25

40. The system of claim 39, wherein, if no other nodes from among the plurality of nodes associated with the load balancer are not known by the load balancer to be inactive, the load balancer is configured to send a message to the higher-level load balancer to disable the load balancer from receiving further requests.

30

41. The system of claim 39, wherein the load balancer is configured to re-load-balancing requests pending with the inactive node among the plurality of nodes associated with the load balancer and not known by the load balancer to be inactive.

5 42. The system of claim 26 wherein the load balancer is configured to send a dummy request to the selected node after selecting the node, and wherein to determine if the selected node is able to service the request, the load balancer is configured to determine if the selected node fails to respond to the dummy request within a timeout period.

10

43. The system of claim 42 wherein the load balancer is configured to mark the selected node as inactive if the selected node fails to respond to the dummy request within the timeout period.

15 44. The system of claim 43, wherein, if no other nodes from among the plurality of nodes associated with the load balancer are not known by the load balancer to be inactive, the load balancer is configured to send a message to the higher-level load balancer to disable the load balancer from receiving further requests.

20 45. The system of claim 43, wherein the load balancer is configured to re-load-balance requests pending with the inactive node among the plurality of nodes associated with the load balancer and not known by the load balancer to be inactive.

25 46. The system of claim 42 wherein the load balancer is configured to send the request to the selected node if the selected node responds to the dummy request within the timeout period.

30 47. The system of claim 46, wherein to determine if the selected node is able to service the request, the load balancer is further configured to determine if the selected node fails to respond to the request within a timeout period.

48. The system of claim 26, wherein to determine if the selected node is able to service the request, the load balancer is configured to receive a message from the selected node indicating that the selected node is disabled.

5

49. The system of claim 48 wherein the load balancer is configured to mark the selected node as inactive upon receiving said message.

50. The system of claim 49 wherein the load balancer is configured to re-load-
10 balance requests pending with the selected node among the plurality of nodes associated with the load balancer and not known by the load balancer to be inactive upon receiving said message.

51. A computer accessible medium, comprising program instructions
15 executable to implement:

a load balancer receiving a request;

the load balancer selecting a node to handle the request from among a plurality of
20 nodes associated with the load balancer and not known by the load balancer to be inactive;

the load balancer determining if the selected node is able to service the request;

25 if the selected node is determined to be unable to service the request, the load balancer selecting another node to handle the request from among the plurality of nodes associated with the load balancer and not known by the load balancer to be inactive.